

An empirical process view of inverse regression

François Portier*

June 2, 2015

Abstract: Most of the methods among the inverse regression literature rely on a slicing of the range of the response variable. Theoretical results are usually shown assuming that (i) the slices are fixed while in practice estimators are constructed with (ii) random slices that contain the same number of observations. In this paper we obtain the asymptotic normality in the case where the slices contains the same number of observations. This issue matter since we find a gap between the asymptotic distributions related to both approaches (i) and (ii). Along this line, we revisit the asymptotic properties of existing methods such as sliced inverse regression and cumulative inverse regression, and we also introduce a bootstrap procedure that reproduce accurately the law of certain Cramér-von Mises test statistics. Our approach is based on the stochastic analysis of some empirical processes that lie close to a certain subspace of interest called the central subspace.

Key words: Dimension reduction; Sliced inverse regression; Cumulative slicing estimation; Weak convergence in $l^\infty(\mathbb{R})$; Bootstrap; Test.

1 Introduction

Dimension reduction is a powerful tool usually employed to synthesise the dependence between two sets of random variables, say (X, Y) where $X \in \mathbb{R}^p$ is called the vector of predictors and $Y \in \mathbb{R}$ is the variable to explain also called the response variable. Dimension reduction can be used to visualize the dependence in high dimensional data [5], as well as to construct accurate estimators of the conditional distribution of Y knowing X [14]. The most common way to model dimension reduction is to assume a certain structure on the conditional distribution of Y given X (see for instance the introduction of [6]). Here we assume that there exists $\beta_0 \in \mathbb{R}^{p \times d_0}$ such that the joint distribution of (X, Y) satisfies

$$P(Y \in A|X) = P(Y \in A|\beta_0^T X), \quad (1)$$

for every Borel set $A \subset \mathbb{R}$. The objective is to estimate the matrix β_0 or rather, because of identifiability reasons [19], the subspace it generates. This subspace is called the central subspace. To this typical semi-parametric problem, many different approaches have been investigated in the past decades [14], [16], [19], [7]. In this paper we follow the idea of inverse regression introduced by Li [19]. In spite of suffering from theoretical restriction on X inverse regression often leads to estimators that are very accurate and computationally efficient. Inverse regression methods are widely spread probably because they provide a reasonable trade-off between accuracy and complexity.

*Institut de statistique, biostatistique et sciences actuarielles, Université catholique de Louvain, Voie du Roman Pays 20, B1348 Louvain-la-Neuve, Belgium. Research supported by Fonds de la Recherche Scientifique (FNRS) A4/5 FC 2779/2014-2017 No. 22342320. Email addresses: francois.portier@uclouvain.be.

Throughout the paper, we will assume that the central subspace is unique. This is known to be true as soon as X has a density ([23], Theorem 1). For more clarity in the statements we introduce the standardized predictors $Z = \Sigma^{-1/2}(X - EX)$ with $\Sigma = \text{var}(X)$. The standardized central subspace, generated by $\Sigma^{1/2}\beta_0$ is denoted by E_c and we let P be the orthogonal projector on E_c .

Inverse regression is based on the following assumption. We say that X satisfies the linearity condition if

$$E(Z|PZ) = PZ, \quad (\text{LC})$$

examples of such distributions include Gaussian distributions, uniform distribution on the sphere, or more generally the class of spherical variables [12]. Li noticed in [19] that under (1) and (LC),

$$E(Z|Y) \in E_c, \quad (2)$$

with probability 1. He then proposed to approximate E_c by estimating the subspace generated by $\text{var}(E(Z|Y))$. The estimation is realized through a slicing of the response Y . A similar slicing method that has been shown to be more efficient is the minimum discrepancy approach (MD) [7]. When facing regression models with a symmetric link function (often refereed as the SIR pathology), SIR is inconsistent. Li suggested in [19] to use second order moments of the predictors. Following this idea, some authors have introduced order 2 moments methods as for instance sliced average variance estimation (SAVE) [8], directional regression [18] and order 2 optimal function [23]. These methods require an additional assumption called the constant covariance condition,

$$\text{var}(Z|PZ) = \text{const.}, \quad (\text{CCV})$$

they are based on the result that, under (1), (LC) and (CCV), it holds that

$$\text{var}(Z|Y) - I \in E_c, \quad (3)$$

where I is the identity matrix.

As a consequence of Equations (2) and (3), the current literature have put the focus on the estimation of subspaces that are generated by conditional quantities. A natural issue which arises is to know whether a nonparametric estimation is really necessary. On the one hand, some authors have studied the limiting distribution of SIR and SAVE estimators as the slicing becomes more thin [17], [30], [29], [21]. Though the conditional quantities $E(Z|Y)$ or $\text{var}(Z|Y)$ can not be estimated at rates root n , these authors shown that the rate root n is in fact available when estimating moments of these quantities, as for instance $\text{var}(E(Z|Y))$. On the other hand, other authors considered a constant number of slices, so that the length of the slices does not go to 0 [8], [7], [18], [23]. In favour of the latter approach, for order 1 moments methods, one might argue that since

$$E[Z\psi(Y)] \in E_c,$$

for any measurable function ψ such that $E[Z\psi(Y)] < \infty$, the whole space E_c will eventually be recovered as soon as the number of function ψ is large (see [23], Theorem 3). Going further, a natural idea is to consider the estimation of E_c when ψ describe a given class of function without necessarily being a slicing. This can be found in [28], where the use of polynomial functions are discussed, and in [23] where the optimal choice of ψ among a Hilbert space is considered (see also [3] for the use of basis functions). In [31], the authors consider a sum over a non-countable

class of functions: the indicators of sets $\{Y \leq t\}$, when t varies on the real line. The underlying method is an integral based method called cumulative slicing estimation (CUME). In this paper we continue along this line by providing an empirical process view of the problem, by indexing the estimators by the elements of a given class of function.

The first contribution of the paper is the introduction and the study of two empirical processes that get closer to E_c as the number of observations increases, one is based on the first conditional moments of Z knowing Y and the other one rely on the second conditional moments of Z knowing Y . Let $\Phi : \mathbb{R} \rightarrow [0, 1]$ be a distribution function and denote by Φ^- its generalized inverse, given by

$$\Phi^-(u) = \inf\{t \in \mathbb{R} : \Phi(t) \geq u\},$$

for every $u \in [0, 1]$. We define the first moment process as

$$c_\Phi(u) = E(Z \mathbb{1}_{\{Y \leq \Phi^-(u)\}}),$$

and the second moment process as

$$C_\Phi(u) = E((ZZ^T - I) \mathbb{1}_{\{Y \leq \Phi^-(u)\}}),$$

for each $u \in [0, 1]$. Clearly, under (1) and (LC), $c_\Phi(u) \in E_c$, if moreover (CCV) holds then $C_\Phi(u) \in E_c$, for every $u \in [0, 1]$.

The fact that c_Φ and C_Φ are indexed by the class of indicator functions plays a key role in our analysis. First the class of indicators is large enough to ensure an exhaustive characterization of E_c . Second it is sufficiently small to enjoy a small metric entropy which is at the root of many nice asymptotic properties of the associated empirical process [27]. Such properties include weak convergence of estimators of c_Φ and C_Φ with root n rates, and the validity of some general weighted bootstrap procedures.

The function Φ is a user-selected function. As in copula modelling, to alleviate the effect of the marginal distribution of Y in the estimation, it is convenient to “uniformize” the variable Y . This is done by choosing Φ equal to F : the cumulative distribution function (cdf) of Y . Since F is unknown, such a choice involves a little more technicalities in the proof but it leads to an accurate and computationally simple rank-based estimator. In [11], the authors studied the weak convergence of the empirical copula process. Following their approach, our theoretical study is based on both the delta-method for stochastic processes and a “trick” allowing us to consider Y as uniformly distributed (see Remark 1).

The study of these processes is conducted in Section 2. It shall be the basis of our study about inverse regression, the main point of which are outlined below.

- i) (see Section 3.1 and 3.2) We obtain the exact asymptotic distribution of SIR when the number of slices is fixed and each slice contains the same number of observations. This way of computing SIR was already pointed out in Remark 4.2 in [19] and it is the most common way to compute slicing estimators. The main issue here is to account for the effect of the randomness of the slices on the asymptotic distribution of SIR. To our knowledge, such results are new in the literature.
- ii) (see Section 3.2) We introduce the class of integral based methods that approximate E_c through the range of the matrices

$$\int \mu(u) \mu(u)^T d\nu(u),$$

where μ stands for a stochastic process that lies in E_c , e.g. c_Ψ or C_Ψ , and ν is a given probability measure. We show that this class includes interesting members such as SIR and CUME. Under mild condition, we prove the asymptotic normality and we provide a valid bootstrap procedure that indeed accounts for the randomness of the slices. Bootstrap is made through a weighting of the estimators that follows from [25], it includes for instance Efron's original bootstrap or the Bayesian bootstrap. Even for SIR or CUME, no such bootstrap was available in the literature.

- iii) (see Section 3.3) In the same spirit as the integral based methods of (ii), we develop several statistical tests of the type Cramér-von Mises: (a) a test of the dimension of a model, i.e. $d_0 = d$ against $d_0 > d$, for some $1 \leq d \leq p$, (b) following [9], a test to assess the no effect of some user-selected sets of predictors, say $\eta^T Z$ with $\eta \in \mathbb{R}^{p \times (p-d)}$, (c) a test similar to (b), but now with η estimated by a given dimension reduction method. The latter might lead us to evaluate whether a model is subject to the SIR pathology. The limiting laws of the considered statistics are fairly hard to estimate so that we provide a valid Bootstrap procedure in order to compute their quantiles. The choice of the bootstrap is crucial in testing since the bootstrap statistic needs to behave similarly as the statistic under H_0 even if H_1 is realized [13]. To implement the bootstrap, we follow ideas from [24] where a constraint bootstrap was developed for testing the rank of a matrix.

A numerical analysis is given in Section 4, in which we study the behaviour of the bootstrap approximation in significance testing.

2 Preliminary results on empirical processes

2.1 Definitions

Using the outer integral, the author Hoffman-Jorgensen has defined a notion of weak convergence of random sequences valued in a metric space [15]. This allows some elements of the considered sequences to be non-measurable provided that their limits are. We equip the space $l^\infty(\mathbb{R})$ of bounded real functions defined on \mathbb{R} with the supremum norm $\|\cdot\|_\infty$. We consider in this paper weak convergence of random elements in $l^\infty(\mathbb{R})$ in the sense of Hoffman-Jorgensen. Let $(Z_i, Y_i)_{1 \leq i \leq n}$ be an i.i.d. sequence of random elements lying in \mathbb{R}^2 with law P . We say that a class of measurable functions $\mathcal{F} \subset l^\infty(\mathbb{R})$ is P -Donsker if

$$n^{-1/2} \sum_{i=1}^n (f(Z_i, Y_i) - Ef(Z_1, Y_1)) \text{ converges weakly in } l^\infty(\mathbb{R}).$$

A complete study of the notion of weak convergence in metric spaces and Donsker classes is proposed in [27]. The following lemma will be useful in the next.

Lemma 1. *Assume that EZ^2 is finite, then $\{(z, y) \mapsto z\mathbb{1}_{(-\infty, t]}(y), t \in \mathbb{R}\}$ is P -Donsker.*

Proof. Let $\mathcal{G} = \{(x, y) \mapsto x\mathbb{1}_{(-\infty, t]}(y), t \in \mathbb{R}\}$. First, it is well-known that $\mathcal{F} = \{\mathbb{1}_{(-\infty, t]}, t \in \mathbb{R}\}$ is P -Donsker (see for instance [27], Example 2.5.4, page 129). In particular, the covering number of \mathcal{F} is such that

$$N(\epsilon, \mathcal{F}, L_2(P)) \leq \frac{2}{\epsilon^2}. \quad (4)$$

Second, since functions of \mathcal{G} have the form $g = \phi(id, f)$ for some $f \in \mathcal{F}$, where $\phi(x, y) = xy$ and id stands for the identity function, we can write

$$\mathcal{G} = \phi(\{id\}, \mathcal{F}).$$

Let f_1 and f_2 be functions in \mathcal{F} , since we have

$$|\phi \circ (id, f_1)(x, y) - \phi \circ (id, f_2)(x, y)|^2 = x^2(f_1(y) - f_2(y))^2,$$

we can apply Theorem 2.10.20 page 199 in [27] (the condition above corresponds to (2.10.19), an envelope for \mathcal{F} is the function equal to 1 everywhere). In view of the bound for the covering number of \mathcal{F} given in (4), and the fact that the covering number of a single element is 1, the uniform entropy condition is checked, making the class \mathcal{G} a P -Donsker class. \square

2.2 Asymptotic behaviour when Φ is known

From now on, $(Z_i, Y_i)_{1 \leq i \leq n}$ is an i.i.d. sequence of random elements lying in $\mathbb{R}^p \times \mathbb{R}$ and drawn from model (1) with $\text{var}(Z_1) = I$ and $EZ_1 = 0$. We denote by $|\cdot|_2$ the Euclidean norm. In what follows, elements of interest belong to the space $l^\infty([0, 1])^p$ that is (with a slight abuse of notation) the space of bounded \mathbb{R}^p -valued functions defined on $[0, 1]$. The empirical processes that estimate the processes c_Φ and C_Φ are defined as follows, for every $u \in [0, 1]$, by

$$\hat{c}_\Phi(u) = \frac{1}{n} \sum_{i=1}^n Z_i \mathbb{1}_{\{Y_i \leq \Phi^-(u)\}} \quad \text{and} \quad \hat{C}_\Phi(u) = \frac{1}{n} \sum_{i=1}^n (Z_i Z_i^T - I) \mathbb{1}_{\{Y_i \leq \Phi^-(u)\}}.$$

We introduce the matrix

$$\gamma_1(u, v) = \text{cov}(Z \mathbb{1}_{\{Y \leq \Phi^-(u)\}}, Z \mathbb{1}_{\{Y \leq \Phi^-(v)\}}).$$

Theorem 2. *Assume that $E[|Z_1|_2^2]$ is finite and Φ is a cdf, then $\sqrt{n}(\hat{c}_\Phi - c_\Phi)$ converges weakly in $l^\infty([0, 1])^p$ to a tight Gaussian process with zero-mean and covariance function γ_1 .*

Proof. Each coordinate of the process $\sqrt{n}(\hat{c}_\Phi - c_\Phi)$ can be written as $\sqrt{n}(\mathbb{P}_n - P)g$ where by Lemma 1, g lies in a Donsker class. Because tightness is equivalent to tightness of each coordinates, it implies that the process $\sqrt{n}(\hat{c}_\Phi - c_\Phi)$ is tight. The limiting process is then given by the limiting distribution of the finite dimensional laws obtained by the multivariate central limit theorem. \square

We now obtain the weak convergence of $\sqrt{n}(\hat{C}_\Phi - C_\Phi)$. To state this we define the operator vec that vectorizes a matrix by stacking its columns, and we introduce the matrix

$$\Gamma_1(u, v) = \text{cov}(\text{vec}(ZZ^T - I) \mathbb{1}_{\{Y \leq \Phi^-(u)\}}, \text{vec}(ZZ^T - I) \mathbb{1}_{\{Y \leq \Phi^-(v)\}}).$$

Corollary 2. *Assume that $E[|Z_1|_2^4]$ is finite and Φ is a cdf, then $\sqrt{n}(\hat{C}_\Phi - C_\Phi)$ converges weakly in $l^\infty([0, 1])^{(p \times p)}$ to a tight Gaussian process with zero-mean and covariance function Γ_1 .*

Proof. We apply Theorem 2 with $\text{vec}(ZZ^T - I)$ in place of Z . \square

2.3 Asymptotic behaviour when Φ is the cdf of Y

We focus on the case where $\Phi = F$ the unknown distribution function of Y . Since

$$c_F(u) = E(Z \mathbb{1}_{\{Y \leq F^-(u)\}}) \quad \text{and} \quad C_F(u) = E((ZZ^T - I) \mathbb{1}_{\{Y \leq F^-(u)\}}),$$

this choice “uniformizes” the variable Y and, as a consequence, vanishes the effect of the distribution of Y on the estimation. Clearly, we can not follow the same path as previously since the estimation of F will certainly affect the limiting process. We introduce the empirical cdf

$$\hat{F}(t) = n^{-1} \sum_{i=1}^n \mathbb{1}_{\{Y_i \leq t\}},$$

defined for each $t \in \mathbb{R}$. Our estimators are plugged-in estimators, i.e. c_F and C_F are respectively estimated by $\hat{c}_{\hat{F}}$ and $\hat{C}_{\hat{F}}$ given by

$$\hat{c}_{\hat{F}}(u) = \frac{1}{n} \sum_{i=1}^n Z_i \mathbb{1}_{\{Y_i \leq \hat{F}^-(u)\}} \quad \text{and} \quad \hat{C}_{\hat{F}}(u) = \frac{1}{n} \sum_{i=1}^n (Z_i Z_i^T - I) \mathbb{1}_{\{Y_i \leq \hat{F}^-(u)\}}.$$

Remark 1. An important point is that when F is continuous, without loss of generality, the variables Y_i 's can be assumed to be uniformly distributed on $[0, 1]$. Equivalently the limiting function \hat{F} can be assumed to be the identity function on $[0, 1]$. To show this, first note that because \hat{F} is a càd-làg function that has $1/n$ -jumps at each Y_i , it is easy to show that for any $u \in [0, 1]$,

$$\{Y_i \leq \hat{F}^-(u)\} \Leftrightarrow \{\hat{F}(Y_i) < u + n^{-1}\}. \quad (5)$$

Then we have that

$$\hat{c}_{\hat{F}}(u) = \frac{1}{n} \sum_{i=1}^n Z_i \mathbb{1}_{\{\hat{F}(Y_i) < u + n^{-1}\}}$$

and a similar expression holds for $\hat{C}_{\hat{F}}$. This makes the previous estimators being sums over the Z_i 's and the rank statistics $\hat{F}(Y_i)$'s. Second note that the rank statistics based on the Y_i 's are equal to the rank statistics based on the uniformized variables $F(Y_i)$'s. As a consequence of this two facts, the processes $\hat{c}_{\hat{F}}$ and $\hat{C}_{\hat{F}}$ can be constructed identically with the samples $(Z_i, Y_i)_{1 \leq i \leq n}$ and $(Z_i, F(Y_i))_{1 \leq i \leq n}$. From now on in the proofs, since $F(Y_i)$ is uniformly distributed on $[0, 1]$ (because of the continuity of F), we can assume without any loss of generality that the variable Y is uniformly distributed.

To compute the asymptotic distribution, since $c_F = c_{id} \circ F^-$, we use the Delta method in metric spaces stated in Theorem 3.9.4 of [27]. This approach has been employed for instance in [27], page 389, and in [11], both in the context of the weak convergence of the empirical copula process. More precisely, we follow this scheme:

- i) Use Lemma 1 to obtain the weak convergence of the process $(s, t) \mapsto n^{1/2}(\hat{F}(s) - F(s), \hat{c}_{id}(t) - c_{id}(t))$.
- ii) Apply the Delta method with the map $(F, c_{id}) \mapsto c_{id} \circ F^-$.

Because the latter map involves the quantile transformation that is not Hadamard differentiable everywhere (see Lemma 3.9.23 in [27]), the fact that \hat{F} can be assumed to converge to the cdf of a uniform distribution (by Remark 1) is a key step in our proof. We define the function $\gamma_2 : [0, 1]^2 \rightarrow \mathbb{R}^{(p+1) \times (p+1)}$ given by

$$\gamma_2(u, v) = \text{cov} \left(\begin{pmatrix} 1 \\ Z \end{pmatrix} \mathbb{1}_{\{Y \leq F^-(u)\}}, \begin{pmatrix} 1 \\ Z \end{pmatrix} \mathbb{1}_{\{Y \leq F^-(v)\}} \right),$$

and $\gamma_3 : [0, 1]^2 \rightarrow \mathbb{R}^{p \times p}$ by

$$\gamma_3(u, v) = (-\partial c_F(u), I) \gamma_2(u, v) (-\partial c_F(v), I)^T,$$

where ∂c_F stands for the derivative of the map $u \mapsto c_F(u)$.

Theorem 3. Assume that $E[|Z_1|_2^2]$ is finite and F is continuous. Then if c_F is continuously differentiable, $\sqrt{n}(\hat{c}_{\hat{F}} - c_F)$ converges weakly in $l^\infty([0, 1])^p$ to a tight Gaussian process with zero-mean and covariance function γ_3 .

Proof. Without loss of generality, we can put $F(Y_i)$ in place of Y_i (see Remark 1). We denote by $id_{[0,1]}$ the cdf of the uniform distribution. By applying Lemma 1, the process $\sqrt{n}(\widehat{G} - G)$, with $\widehat{G}(s, t) = (\widehat{F}(s), \widehat{c}_F(t))$ and $G(s, t) = (id_{[0,1]}(s), c_F(t))$, converges weakly in $l^\infty(\mathbb{R}) \times l^\infty([0, 1])^p$ to a tight Gaussian element. Now since

$$\widehat{c}_{\widehat{F}} = \psi(\widehat{G}) \quad \text{and} \quad c_F = \psi(G), \quad (6)$$

where $\psi : \mathcal{F} \times l^\infty([0, 1])^p \rightarrow l^\infty([0, 1])^p$, \mathcal{F} being the space of cdf with support included in $[0, 1]$, is given by

$$\psi : (f_1, f_2) \mapsto (f_1^-, f_2) \mapsto f_2 \circ f_1^-, \quad (7)$$

we can apply Theorem 3.9.4, page 374 in [27] which basically says that $\sqrt{n}(\psi(\widehat{G}) - \psi(G))$ is P -Donsker provided that the map ψ is Hadamard differentiable. In what follows, we first show that ψ is Hadamard differentiable, and then we compute the asymptotic variance. Using Lemma 3.9.23, assertion (ii), page 386 in [27], the first map of Equation (7) reduced to $f \mapsto f^-$ is Hadamard differentiable at the function $id_{[0,1]}$ tangentially to $C[0, 1]$. Moreover its derivative at $id_{[0,1]}$, in the direction h_1 is given by $-h_1$. Since c_F is Fréchet differentiable, by Lemma 3.9.27, page 388 in [27], the second map in Equation (7) is Hadamard differentiable at $(id_{[0,1]}^-, c_F)$, tangentially to $C[0, 1]$ (because continuous functions are uniformly continuous on compacts). Its derivative at $(id_{[0,1]}^-, c_F)$, in the direction (h_1, h_2) , is given by $h_1 \times \partial c_F + h_2$. By the chain rule, the function ψ is Hadamard differentiable at the point $(id_{[0,1]}, c_F)$ tangentially to $C[0, 1]$. At this point, in the direction (h_1, h_2) , its derivative is given by $-h_1 \times \partial c_F + h_2$. Hence, the limiting process has the representation

$$u \mapsto w_1(u) - \partial c_F(u) \times B(u) = (-\partial c_F(u), I) \begin{pmatrix} B(u) \\ w_1(u) \end{pmatrix},$$

where (B, w_1) is the Gaussian limit of $\widehat{H} : u \mapsto \sqrt{n}(\widehat{G} - G) \circ (u, u)$. Its covariance function is computed by applying the central limit theorem that gives

$$(\widehat{H}(u_1), \dots, \widehat{H}(u_K)) \xrightarrow{d} ((B, w_1)(u_1), \dots, (B, w_1)(u_K)),$$

where $\text{vec}((B, w_1)(u_1), \dots, (B, w_1)(u_K))$ is a Gaussian vector with mean 0 and covariance matrix having the block decomposition $(\gamma_2(u_k, u_l))_{1 \leq k, l \leq K}$. □

To obtain a similar result about the order 2 moments process, we define the function $\Gamma_2 : [0, 1]^2 \rightarrow \mathbb{R}^{(p+1) \times (p+1)}$ by

$$\Gamma_2(u, v) = \text{cov} \left(\begin{pmatrix} 1 \\ \text{vec}(ZZ^T - I) \end{pmatrix} \mathbb{1}_{\{Y \leq F^{-1}(u)\}}, \begin{pmatrix} 1 \\ \text{vec}(ZZ^T - I) \end{pmatrix} \mathbb{1}_{\{Y \leq F^{-1}(v)\}} \right),$$

and $\Gamma_3 : [0, 1]^2 \rightarrow \mathbb{R}^{p \times p}$ by

$$\Gamma_3(u, v) = (-\partial \text{vec}(C_F)(u), I) \Gamma_2(u, v) (-\partial \text{vec}(C_F)(v), I)^T.$$

where $\partial \text{vec}(C_F)(u)$ stands for the derivative of the map $u \mapsto \text{vec}(C_F)(u)$.

Corollary 3. *Assume that $E[|Z_1|^4]$ is finite and F is continuous. Then if $\text{vec}(C_F)$ is continuously differentiable, $\sqrt{n}(\widehat{C}_{\widehat{F}} - C_F)$ converges weakly in $l^\infty([0, 1])^{(p \times p)}$ to a tight Gaussian process with zero-mean and covariance function Γ_3 .*

Proof. We apply Theorem 2 with $\text{vec}(ZZ^T - I)$ in place of Z . □

2.4 The Bootstrap

In light of the limiting covariance processes given in the previous section, in particular because of the presence of ∂c_F and $\partial \text{vec}(C_F)$ but also the possibly high-dimensionality of these processes, the asymptotic distributions are fairly hard to estimate. As a consequence, for making inference, it seems necessary to develop a bootstrap strategy. Efron [10] introduced the original bootstrap that consists in a sampling with equi-probability and replacement of the original sample. In [25], the authors considered a more general re-sampling plan based on weights $w_{i,n}$, $i = 1, \dots, n$ that verified

(B1) the random sequence $(w_{i,n})_{1 \leq i \leq n}$ is exchangeable, i.e. for every permutation (π_1, \dots, π_n) of $(1, \dots, n)$, $(w_{i,n})_{1 \leq i \leq n}$ has the same law as $(w_{\pi_i,n})_{1 \leq i \leq n}$,

(B2) denote by S_n the survival function of $w_{1,n}$, we have

$$\sup_{n \geq 1} \int S_n(u)^{1/2} du < +\infty \quad \text{and} \quad \lim_{A \rightarrow +\infty} \limsup_{n \rightarrow +\infty} \sup_{t \geq A} t^2 S_n(t) = 0.$$

(B3) $w_{i,n} \geq 0$, $\sum_{i=1}^n w_{i,n} = n$, $n^{-1} \sum_{i=1}^n (w_{i,n} - 1)^2 \xrightarrow{P} 1$.

Examples of such weights, are given in [25]. Now we define the bootstrap processes

$$\begin{aligned} \hat{F}^*(t) &= n^{-1} \sum_{i=1}^n w_{i,n} \mathbb{1}_{\{Y_i \leq t\}}, \\ \hat{c}_\Phi^*(u) &= \frac{1}{n} \sum_{i=1}^n w_{i,n} Z_i \mathbb{1}_{\{Y_i \leq \Phi^-(u)\}}, \end{aligned}$$

for every $t \in \mathbb{R}$ and every $u \in [0, 1]$. The bootstrap of \hat{c}_Φ (resp. $\hat{c}_{\hat{F}}$) is made by \hat{c}_Φ^* (resp. $\hat{c}_{\hat{F}}^*$). The following theorem basically says that the bootstrap in probability (in the sense of [25]) works.

Theorem 4. *Under (B1) to (B3), assume that $E[|Z_1|^2]$ is finite and Φ is a cdf, then conditionally on the sample,*

$n^{1/2}(\hat{c}_\Phi^ - \hat{c}_\Phi)$ has the same weak limit as $n^{1/2}(\hat{c}_\Phi - c_\Phi)$, in probability.*

If moreover F is continuous and c_F is continuously differentiable, then conditionally on the sample,

$n^{1/2}(\hat{c}_{\hat{F}}^ - \hat{c}_{\hat{F}})$ has the same weak limit as $n^{1/2}(\hat{c}_{\hat{F}} - c_F)$, in probability.*

Proof. The first statement is a direct consequence of Lemma 1 and Theorem 2.1 in [25]. For the second statement, we first apply the trick detailed in Remark 1 to the bootstrap estimator. Indeed it is easy to see that $\hat{c}_{\hat{F}}^*$ can be constructed as well from the sample $(Z_i, F(Y_i))_{1 \leq i \leq n}$, so that the weak limit of \hat{F}^* can be assumed to be $id_{[0,1]}$. This is due to the equivalence between

$$\{Y_i \leq \hat{F}^{*-}(u)\} \Leftrightarrow \{\hat{F}^*(Y_i) < u + n^{-1}w_{i,n}\},$$

for any $u \in [0, 1]$, plus the fact that the bootstrap ranks $\hat{F}^*(Y_i)$'s are the same as the uniformized bootstrap ranks (i.e. based on the $F(Y_i)$'s rather than the Y_i 's). Then by applying again Lemma 1 with Theorem 2.1 in Paestgrad and Wellner, the process $\sqrt{n}(\hat{G}^* - \hat{G})$, with $\hat{G}^*(s, t) = (\hat{F}^*(s), \hat{c}_{\hat{F}}^*(t))$, has the same limiting distribution as $\sqrt{n}(\hat{G} - G)$ (defined in the proof of Theorem 3), that is a tight Gaussian element of $l^\infty(\mathbb{R}) \times l^\infty([0, 1])^p$. Then we can invoke the Delta-method for the bootstrap stated as Theorem 3.9.11, page 378, in [27]. \square

Similarly, we define \hat{C}_Φ^* by

$$\hat{C}_\Phi^*(u) = \frac{1}{n} \sum_{i=1}^n w_{i,n} (Z_i Z_i^T - I) \mathbb{1}_{\{Y_i \leq \Phi^-(u)\}},$$

for every $u \in [0, 1]$, and we obtain this corollary.

Corollary 4. *Under (B1) to (B3), assume that $E[|Z_1|_2^4]$ is finite and Φ is a cdf, then conditionally on the sample,*

$\sqrt{n}(\hat{C}_\Phi^ - \hat{C}_\Phi)$ has the same weak limit as $\sqrt{n}(\hat{C}_\Phi - C_\Phi)$, in probability.*

If moreover F is continuous and $\text{vec}(C_F)$ is continuously differentiable, then conditionally on the sample,

$\sqrt{n}(\hat{C}_{\hat{F}}^ - \hat{C}_{\hat{F}})$ has the same weak limit as $\sqrt{n}(\hat{C}_{\hat{F}} - C_F)$, in probability.*

Proof. We apply Theorem 2 with $\text{vec}(ZZ^T - I)$ in place of Z . □

3 Application to inverse regression

In this section we are based on the results of the previous section in order to (i) raise some new points about the asymptotics of SIR, (ii) develop a unified framework for inverse regression and (iii) study new bootstrap testing procedure. The variables Z_i 's are assumed to be standardized in order to clarify the statements of the results. In practice we must account for the error induced by estimations of the mean and the variance (see Section 4 for more details).

3.1 Revisiting sliced inverse regression

Sliced inverse regression [19] is based on the vectors

$$n^{-1} \sum_{i=1}^n Z_i \mathbb{1}_{\{Y_i \in I(h)\}},$$

where $I(h)$, for $h = 1, \dots, H$ is a partition of the range of the Y_i 's. In practice, to diminish the chance of having a poor estimation of such vectors, it is convenient to keep the same number of observations within each slice (this was already pointed-out in Remark 4.2 by [19] and this is how SIR is usually run). Consequently each member $I(h)$ of the partition is random because it depends on the Y_i 's. Meanwhile when describing the asymptotic behaviour, many authors have ignored this additional source of randomness (see among others [9], [7] or [23]). In what follows, we show that the randomness of the partition $I(h)$ can not be neglected since we find that it participates in the asymptotic variance of the estimation. Our approach can work because the slicing $I(h)$ is expressed in a simple way with the help of the rank statistics $\hat{F}(Y_i)$'s. Hence we shall apply in the next, Theorem 3 and Corollary 3. For brevity, we focus on SIR, but the same analysis can be extended to second order slicing methods such as for instance, SAVE and DR.

Consider a multi-slice procedure with H slices. Denote by $\lceil \alpha \rceil$ the smallest integer greater than or equal to α . A reasonable way to dispatch the data among the slices should be with $\lceil n/H \rceil$ observations in the first slice, $\lceil 2n/H \rceil - \lceil n/H \rceil$ in the second, ..., $n - \lceil n(H-1)/H \rceil$ in the last slice. Note that as soon as n is a multiple of H , each slice contains exactly the same number of observations n/H . The SIR estimator is the subspace generated by

$$(\hat{c}_{\hat{F}}(u_1), \hat{c}_{\hat{F}}(u_2) - \hat{c}_{\hat{F}}(u_1), \dots, \hat{c}_{\hat{F}}(u_H) - \hat{c}_{\hat{F}}(u_{H-1})),$$

with $u_h = h/H$, and the corresponding estimation with nonrandom slices is the span of the matrix

$$(\widehat{c}_F(u_1), \widehat{c}_F(u_2) - \widehat{c}_F(u_1), \dots, \widehat{c}_F(u_H) - \widehat{c}_F(u_{H-1})).$$

Invoking Theorems 2 and 3, for any $h \in \{1, \dots, H\}$, because the sequences $\widehat{c}_{\widehat{F}}(u_h)$ and $\widehat{c}_F(u_h)$ have a different asymptotic distribution, the latter matrices neither. To highlight differences in the behaviour of $\widehat{c}_{\widehat{F}}$ and \widehat{c}_F , we consider the following tool model

$$Y = X_1 + .1e, \tag{8}$$

where $(X, e) \in \mathbb{R}^5$ follows a standard normal distribution. In order to keep clear our statements and conclusions, we focus on the first slice of SIR in which the number of observations $\lceil nu \rceil$ varies with u from 1/2 to 0. There are two different ways to compute it:

- Order the responses Y_i 's, create a slice containing the first $\lceil nu \rceil$ observations, compute the mean over the X_i 's within the slice. This gives the vector $\widehat{c}_1 = \widehat{c}_{\widehat{F}}(u)$.
- Create a slice according to $Y_i \leq F^-(u)$ (the slice is independent of the observations Y_i 's), compute the mean over the X_i 's within the slice. This gives the vector $\widehat{c}_2 = \widehat{c}_F(u)$.

By means of simulations, we evaluate the first coordinate of the latter quantities 1000 times. The resulting boxplots, for different values of n are reported in Figure 1.

Starting from $u = 1/2$ (meaning that the observations have been cut in half), where both variances are the same, we see that, as u decreases, the dispersion of $\widehat{c}_2(u)$ becomes larger, whereas it is clearly more stable for $\widehat{c}_1(u)$. Note also that whereas $\widehat{c}_2(u)$ is unbiased, $\widehat{c}_1(u)$ suffers from a slight bias in small sample sizes. This sheds light on two things: (i) the limiting distribution of $\widehat{c}_1(u)$ and $\widehat{c}_2(u)$ are different and one should care about that as soon as inference is of matter, (ii) for this generic example, $\widehat{c}_1(u)$ is more efficient than $\widehat{c}_2(u)$, highlighting that it is more accurate to have a control on the number of observations within the slices. To the best of our knowledge, the latter question is still open although this is only theoretical because the matrix based on c_2 is not even computable (unless we know the law of Y).

Remark 2. In light of Theorems 2 and 3, it is the function ∂c_F that determines whether the asymptotic is affected by the randomness of the slices. In the case of an additive regression model $Y = g(\beta_0^T X) + e$ with $e \perp\!\!\!\perp X$, we find that

$$c_F(u) = E[Z F_e(F^-(u) - g(\beta_0^T X))],$$

where F_e is the cdf of e . At $u = 0$ and $u = 1$ this quantity equals 0, then under the assumption of Theorem 3, by the Rolle's theorem there exists at least one $v \in (0, 1)$ such that $\partial c_F(v) = 0$. As a consequence, the asymptotic distributions of $\widehat{c}_{\widehat{F}}(v)$ and $\widehat{c}_F(v)$ are the same. Nevertheless this certainly will not happen at each slice boundary u_h , as it is highlighted in Figure 1 for Model (8), for which $v = 1/2$.

Remark 3 (cumulative slicing estimation). In [31], the authors consider spaces generated by the integral

$$\int_0^1 \widehat{c}_{\widehat{F}}(u) \widehat{c}_{\widehat{F}}(u)^T du.$$

Contrary to most of the existing methods, any slicing is no longer necessary. They focus on a large p small n context and give a limit theorem by borrowing a U -statistic approach, that is rather different than our empirical process approach. Their simulation results highlight that CUME is competitive with SIR and performs even better in several situations.

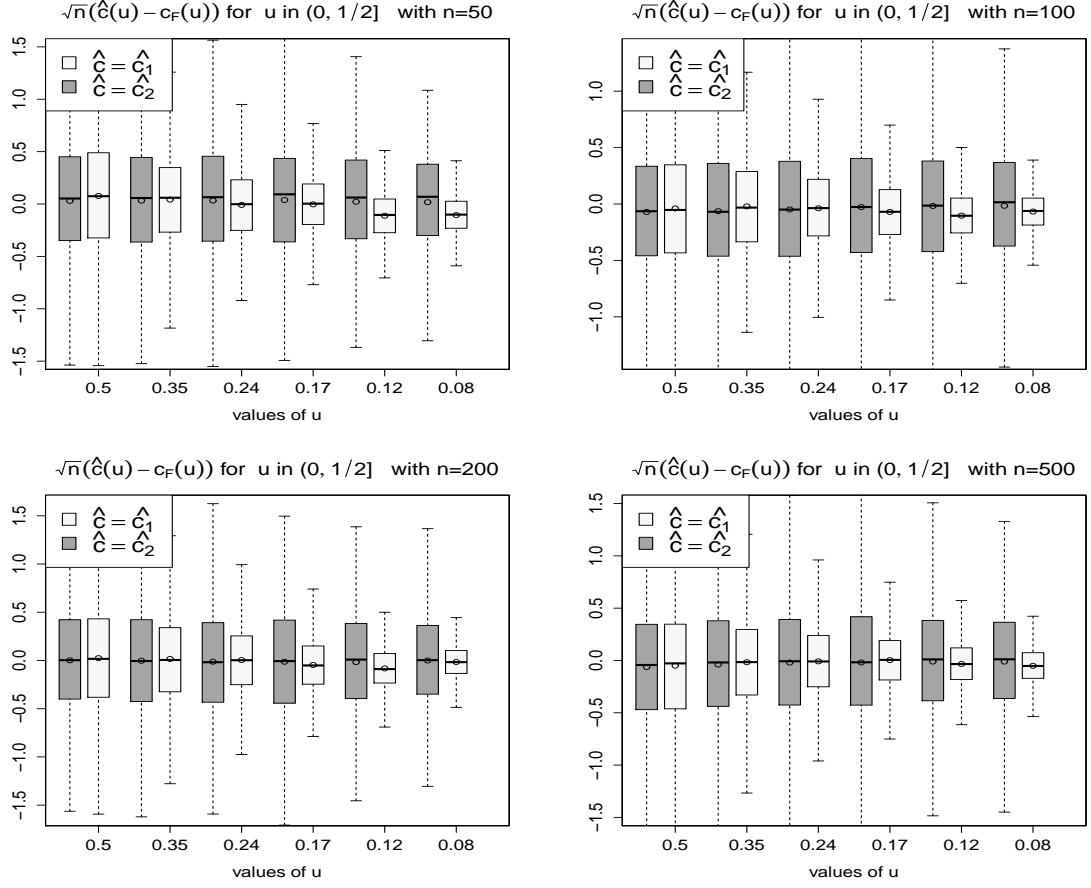


Figure 1: Boxplots of the law of the first coordinate of $\sqrt{n}(\hat{c}_k(u) - c_F(u))$, for $k = 1, 2$, based on 1000 replications.

3.2 Integral approach: a unified framework

In order to consider in the next a broad class of different methods, and in particular to include SIR and CUME, it is useful to introduce the matrix

$$A_\nu(\hat{\mu}) = \int \hat{\mu}(u) \hat{\mu}(u)^T d\nu(u),$$

where $\hat{\mu}$ belongs to the space $l^\infty([0, 1])^{(p \times q)}$ with $q \geq 1$ and ν is a probability measure on $[0, 1]$. In order to estimate E_c , the process $\hat{\mu}$ shall be a combination of processes studied in the Section 2, namely $\hat{c}_{\hat{F}}$ and $\hat{C}_{\hat{F}}$. As soon as (LC) and (CCV) are realized, the corresponding limit of $A_\nu(\hat{\mu})$ generates a subspace of E_c . Hence the estimation of E_c follows from an eigendecomposition of $A_\nu(\hat{\mu})$ by taking the eigenvectors associated to the d_0 largest eigenvalues as the estimated basis of E_c .

Order 1 moments based methods. One easily sees that taking $\hat{\mu}$ equal to $\hat{c}_{\hat{F}}(u)$ and ν equal to the uniform distribution on $[0, 1]$, leads to CUME. Now let $\pi \in [0, 1]$ and define the quantities

$$\begin{aligned} \hat{m}(u, \pi) &= \hat{c}_{\hat{F}}(u + \pi/2) - \hat{c}_{\hat{F}}(u - \pi/2), \\ m(u, \pi) &= c_F(u + \pi/2) - c_F(u - \pi/2), \end{aligned}$$

the SIR estimators can be expressed as the space generated by

$$\sum_{h=1}^H \hat{m}((2h-1)/2H, 1/H) \hat{m}((2h-1)/(2H), 1/H)^T.$$

This is a direct consequence of the definition of SIR given in the previous section. As a result, SIR with H slices belongs to our framework. It corresponds to the matrix $\hat{A}_\nu(\hat{\mu})$ when $\hat{\mu}$ equal to $\hat{m}(\cdot, 1/H)$ and ν is the cdf of a discrete uniform random variable over the set $\{(2h-1)/(2H), h = 1, \dots, H\}$. Our framework permits also a slight modification of SIR, based on the same process $\hat{m}(\cdot, \pi)$ but with ν being the cdf of a continuous (rather than discrete) uniform random variable on $[0, 1]$.

Order 2 moments based methods. As it is well-known in the literature, SIR and CUME are inconsistent in estimating directions that present a symmetric relationship with the variable Y . To remedy this problem, one can rather consider order 2 moments of the predictor as in SAVE or DR. Within our framework, this means computing $A_\nu(\hat{\mu})$ with $\hat{\mu}$ equal to $\hat{C}_{\hat{F}}$ (cumulative version) or $\hat{M}(u, \pi) = \hat{C}_{\hat{F}}(u + \pi/2) - \hat{C}_{\hat{F}}(u - \pi/2)$ (slicing version).

Denoting by μ the limit in probability of $\hat{\mu}$, we obtain the weak convergence of $n^{1/2}(A_\nu(\hat{\mu}) - A_\nu(\mu))$ by following these steps:

- (A) Weak convergence of the process $n^{1/2}(\hat{\mu} - \mu)$ in $l^\infty([0, 1])^{(p \times q)}$.
- (B) Application of the continuous mapping theorem to extend the convergence to some integral maps.

In Section 2, we have focused on the first step, so that the proof of the following theorem essentially consists in showing the second step. For brevity we formally state our results for the order 1 moments based methods, the extension to order 2 moments based methods being straightforward (see Remark 6).

Theorem 5. Assume that $E[|Z_1|^2]$ is finite, F is continuous and c_F is continuously differentiable, then

(i) if $\hat{\mu} = \hat{c}_{\hat{F}}$, $\mu = c_F$,

$$n^{1/2}(A_\nu(\hat{\mu}) - A_\nu(\mu)) \quad \text{has Gaussian limit} \quad \int w_3(u) c_F^T + c_F w_3(u)^T d\nu(u),$$

(ii) if $\pi \in [0, 1]$, $\hat{\mu} = \hat{m}(\cdot, \pi)$, $\mu = m(\cdot, \pi)$,

$$n^{1/2}(A_\nu(\hat{\mu}) - A_\nu(\mu)) \quad \text{has Gaussian limit} \quad \int w_4(u) m(u, \pi)^T + m(u, \pi) w_4(u)^T d\nu(u),$$

where $w_4(u) = w_3(u + \pi/2) - w_3(u - \pi/2)$ and w_3 is a Gaussian process with covariance γ_3 .

Proof. Since the proofs of (i) and (ii) are very similar we focus on (ii). Invoking Theorem 2 and the continuous mapping theorem stated for instance in [27], page 20, as Theorem 1.3.6, we obtain the weak convergence of $\{\sqrt{n}(\hat{m}(u, \pi) - m(u, \pi))\}_{u \in [0, 1]}$ to the Gaussian process w_4 . For every $u \in [0, 1]$ and $\pi \in [0, 1]$, one can write

$$\begin{aligned} & \hat{m}(u, \pi) \hat{m}(u, \pi)^T - m(u, \pi) m(u, \pi)^T \\ &= (\hat{m}(u, \pi) - m(u, \pi)) m(u, \pi)^T + m(u, \pi) (\hat{m}(u, \pi) - m(u, \pi))^T \\ & \quad + (\hat{m}(u, \pi) - m(u, \pi)) (\hat{m}(u, \pi) - m(u, \pi))^T, \end{aligned}$$

then, as a consequence of the Delta-method, $\{\sqrt{n}(\hat{m}(u, \pi) \hat{m}(u, \pi)^T - m(u, \pi) m(u, \pi)^T)\}_{u \in [0, 1]}$ converges weakly to $\{w_4(u) m(u, \pi)^T + m(u, \pi) w_4(u)^T\}_{u \in [0, 1]}$. Finally applying the continuous mapping theorem to the previous process with the map $f \mapsto \int f(u) d\nu(u)$, we obtain the statement of the theorem. \square

Remark 4 (coverage property). A comparison between the spaces generated by CUME and SIR is relevant to highlight the differences between continuous and discrete methods. The space that SIR estimates is

$$E_{\text{SIR}}^H = \text{span}\{m((2h-1)/2H, 1/H), \quad h = 1, \dots, H\},$$

and under the conditions of Theorem 3 in [23], for H sufficiently large, $E_{\text{SIR}}^H = E_c$. This result is important because it ensures that when H increases, SIR eventually estimates the whole subspace. Nevertheless, this is not sufficient to guarantee a complete estimation of E_c since in practice, we do not know how to choose H . The space estimated by CUME is

$$E_{\text{CUME}} = \text{span}\{c_F(u), \quad u \in [0, 1]\}.$$

It follows that $E_{\text{SIR}}^H \subset E_{\text{CUME}} \subset E_c$. As a consequence, compared with SIR, the method CUME is more likely to recover a larger subspace within E_c .

The bootstrap is made through

$$A_\nu(\hat{\mu}^*) = \int \hat{\mu}^*(u) \hat{\mu}^{*T}(u) d\nu(u),$$

where $\hat{\mu}^*$ is a bootstrap version of μ that can be chosen according to the next theorem. We define the bootstrap process $\hat{m}^*(u, \pi) = \hat{c}_{\hat{F}^*}^*(u + \pi/2) - \hat{c}_{\hat{F}^*}^*(u - \pi/2)$.

Theorem 6. *Under (B1) to (B3), assume that F is continuous and c_F is continuously differentiable, then conditionally on the sample,*

$$n^{1/2}(A_\nu(\hat{\mu}^*) - A_\nu(\hat{\mu})) \quad \text{has the same weak limit as} \quad n^{1/2}(A_\nu(\hat{\mu}) - A_\nu(\mu)), \text{ in probability,}$$

provided that (i) $\hat{\mu}^ = \hat{c}_{\hat{F}^*}^*$ and $\hat{\mu} = \hat{c}_{\hat{F}}$ or (ii) $\hat{\mu}^* = \hat{m}^*(\cdot, \pi)$ and $\hat{\mu} = \hat{m}(\cdot, \pi)$.*

Proof. The proof is similar as the proof of Theorem 5 with the following changes: consider the probability space conditional on the (Y_i, Z_i) 's and replace $\hat{m}(u, \pi)$ by $\hat{m}^*(u, \pi)$ and $m(u, \pi)$ by $\hat{m}(u, \pi)$. \square

Remark 5 (bootstrapping the slices). An accurate description of the asymptotic distribution is necessary for making precise the inference. By Theorem 6, our bootstrap procedure is valid and therefore, shall be use to make inference on $A_\nu(\hat{\mu})$. This is mainly due to the fact that the randomness of the slices has been reproduced by bootstrapping also the estimated cdf of Y , e.g. contrary to $\hat{c}_{\hat{F}^*}^*$, the process $\hat{c}_{\hat{F}}^*$ won't produce a valid bootstrap. Hence, other bootstrap techniques that ignore this randomness will fail in bootstrapping the law of $A_\nu(\hat{\mu})$. Existing bootstrap methods for SIR ([1], [23]) consider the slices as fix, and so they are unable to reproduce correctly the law of SIR as it is usually computed. Nevertheless, when testing specific properties of E_c , it could happen that both bootstrap, respectively directed by $\hat{c}_{\hat{F}^*}^*$ and $\hat{c}_{\hat{F}}^*$, work (see Section 3.3 for more details).

Remark 6 (order 2 moments based methods). Assuming that $E[|Z_1|^4]$ is finite, F is continuous and $\text{vec}(C_F)$ is continuously differentiable, it is an easy exercise to obtain a similar statement as in Theorem 5 (replacing w_3 by W_3 , invoking Corollary 3 rather than Theorem 3) as well as in Theorem 6 (replacing c by C , invoking Corollary 4 rather than Theorem 4).

3.3 Cramér-von Mises tests

The integral methods of the previous section, such as SIR and CUME, produce accurate estimations of E_c (see for instance the simulation study in [31]). Nevertheless, the asymptotic distribution of these methods was unknown from the researchers making difficult any inference based on the matrix $A_\nu(\hat{\mu})$. On the one hand, some authors neglected the effect of the randomness of the slices for SIR ([9], [7] or [23]), on the other hand, other ones employed in addition the Bentler and Xie's approximation [2] in order to compute the asymptotic distribution (see [4] and [24]). Here based on the empirical process approach of Section 2, the purpose is to demonstrate rigorously that bootstrap leads to accurate inference when testing structural properties of E_c . We introduce three tests that asses: the dimension of E_c , the no effect of a set of predictors and the contribution of a given method. At the end of the section, we show that all the tests considered are consistent and that bootstrap is valid to compute their quantiles. All the test statistics that we introduce are of the Cramér-von Mises type, i.e. of the form

$$\int |\hat{f}(u)|_F^2 d\nu(u),$$

where \hat{f} is a certain process that belongs to $l^\infty([0, 1])^{(p \times q)}$ and $|\cdot|_F$ is the Frobenius norm. In our precise situation, because the integrands are piecewise constant, closed-formulas are available, making the tests computationally feasible. This generally no longer happen for Kolmogorov type statistics.

3.3.1 Testing dimensionality

In order to determine the dimension d_0 of E_c , it is usual to test whether d_0 equals a given number, say d , against the alternative d_0 is larger than d , i.e.

$$H_0 : d_0 = d \quad \text{against} \quad H_1 : d_0 > d. \quad (9)$$

Then starting with $d = 0$, if rejected we put $d := d + 1$, until the first acceptance. Different approaches that could be use are summarized in [4] and [24]. In the following we focus on the most common test statistic, based on the sum of eigenvalues of $A_\nu(\hat{\mu})$, given by

$$\hat{\Lambda}_1 = n \sum_{k=d+1}^p \hat{\lambda}_k,$$

where the $\hat{\lambda}_k$'s are the eigenvalues of the matrix $A_\nu(\hat{\mu})$, arranged in decreasing order. We have the formula

$$\hat{\Lambda}_1 = n \text{trace}(\hat{Q} A_\nu(\hat{\mu}) \hat{Q}) = n \int |\hat{Q} \hat{\mu}|_F^2 d\nu(u),$$

where \hat{Q} is the eigenprojector on the eigenspace associated to the $p - d$ smallest eigenvalues of $A_\nu(\hat{\mu})$.

3.3.2 Testing a predictor contribution

Following [9], we develop tests of no effect, on the response variable Y , of a selected group of predictor, say $\eta^T Z$ where $\eta \in \mathbb{R}^{p \times (p-d)}$ is such that $\eta^T \eta = I$. We define β such that $(\beta, \eta) \in \mathbb{R}^{p \times p}$ is an orthogonal matrix. We say that $\eta^T Z$ has no effect on Y if

$$P(Y \in A | \beta^T Z) = P(Y \in A | Z),$$

for any Borel set A . By [9], Proposition 1, this is equivalent to $\eta \in E_c^\perp$. As a consequence, we introduce the hypotheses

$$H_0 : \eta \in E_c^\perp \quad \text{against} \quad H_1 : \eta \notin E_c^\perp. \quad (10)$$

Under the so-called coverage condition, that basically says that E_c is spanned by $A_\nu(\mu)$, the previous set of hypotheses is equivalent to

$$H_0 : \eta^T A_\nu(\mu) \eta = 0 \quad \text{against} \quad H_1 : \eta^T A_\nu(\mu) \eta \neq 0.$$

Therefore a natural statistic for testing H_0 is

$$\hat{\Lambda}_2 = n \text{trace}(\eta^T A_\nu(\hat{\mu}) \eta) = n \int |\eta^T \hat{\mu}(u)|_F^2 d\nu(u).$$

3.3.3 Testing a method contribution

Here we consider a given method whose estimated basis is noted $\hat{\beta} \in \mathbb{R}^{p \times d}$. Let us assume that there exists a basis $\beta \in \mathbb{R}^{p \times d}$ such that $P_{\hat{\beta}}$ converges in probability to P_β , with the notation $P_\beta = \beta \beta^T$. We want to test whether the method misses a direction (asymptotically), i.e.

$$H_0 : \eta \in E_c^\perp \quad \text{against} \quad H_1 : \eta \notin E_c^\perp, \quad (11)$$

where $(\beta, \eta) \in \mathbb{R}^{p \times p}$ is an orthogonal matrix. Let $\hat{\eta}$ be such that $(\hat{\beta}, \hat{\eta})$ is an orthogonal matrix, our statistic is given by

$$\hat{\Lambda}_3 = n \text{trace}(\hat{\eta}^T A_\nu(\hat{\mu}) \hat{\eta}) = n \int |\hat{\eta}^T \hat{\mu}(u)|_F^2 d\nu(u).$$

We have in mind two typical applications. First we aim at testing the so called SIR pathology, i.e. whether an order 1 moments based method fails in recovering the whole subspace. For that purpose, $\hat{\beta}$ might be for instance the estimated basis of SIR or CUME and $\hat{\mu}$ should be based on the order 2 process \hat{C} . Clearly if the model is subject to the order 1 pathology, the test shall reject H_0 . Second the latter procedure can be applied to select the estimated directions for the order 2 optimal function method introduced in [23]. This method alleviates the assumption CCV and produces accurate estimates but a classical eigenvalue-based selection of the directions fails. The initial test of independence developed in [23] rely on a null hypothesis that is too strong. It is more accurate to apply the above test when $\hat{\beta}$ is the estimated basis of the order 2 optimal function method.

3.3.4 Consistency of the tests

Theoretically, a test is said to be consistent if, as n increase, the level converges to the nominal level and the power goes to 1. As it will be stressed out, every of the tests considered previously is consistent. Practically one needs to compute the quantiles of the asymptotic law of the statistic. In our case, those quantiles are difficult to estimate and this could diminish the accuracy of the test [24]. As a consequence, we recommend a bootstrap strategy for computing these quantiles and we show in the next the consistency of our bootstrap procedure.

For the sake of generality, we study all the tests (9), (10) and (11) introduced in the previous section. The statistics $\hat{\Lambda}_k$ for $k = 1, 2, 3$, can be written as follows

$$\hat{\Lambda}_k = n \int |\hat{Q}_k \hat{\mu}(u)|_F^2 d\nu(u),$$

with \hat{Q}_1 the eigenprojector associated to the $p - d$ smallest eigenvalues of $\int \hat{\mu}(u) \hat{\mu}(u)^T d\nu(u)$, $\hat{Q}_2 = \eta \eta^T$ with $\eta \in \mathbb{R}^{p \times (p-d)}$ a basis, and $\hat{Q}_3 = \hat{\eta} \hat{\eta}^T$ the orthogonal projector on the orthogonal complement of the estimated space of a given method as it is described in Section 3.3.3. We also introduce (when they exist) Q_1 , the eigenprojector associated to the $p - d$ smallest eigenvalues of $\int \mu(u) \mu(u)^T d\nu(u)$, $Q_2 = \eta \eta^T$, and Q_3 , the limit of \hat{Q}_3 .

Bootstrap testing requires particular care so that the bootstrap estimator mimics the hypothesis H_0 even when H_1 is realized [13], [24]. For statistics of a similar type as $\hat{\Lambda}_1$, [24] shows that the quantiles can be computed using the technique of the constraint bootstrap. Following their approach, we define the bootstrap statistics $\hat{\Lambda}_k^*$'s by

$$\hat{\Lambda}_k^* = n \int |\hat{Q}_k^* \hat{\mu}_k^*(u)|_F^2 d\nu(u) \quad \text{for } k = 1, 2, 3,$$

with \hat{Q}_1^* the eigenprojector associated to the $p - d$ smallest eigenvalues of $\int \hat{\mu}_1^*(u) \hat{\mu}_1^*(u)^T d\nu(u)$, $\hat{Q}_2^* = \eta \eta^T$, \hat{Q}_3^* is a bootstrap version of \hat{Q}_3 , and for every $u \in [0, 1]$,

$$\hat{\mu}_k^*(u) = (I - \hat{Q}_k) \hat{\mu}(u) + (\hat{\mu}^*(u) - \hat{\mu}(u)). \quad (12)$$

The latter formula is the cornerstone of the bootstrap procedure. It ensures that the bootstrap process $\hat{\mu}_k^*$ is asymptotically contained in a subspace of dimension d , making the bootstrap process having a H_0 -likely behaviour. To guarantee the consistency of the tests, we introduce the following assumptions. A discussion is postponed latter.

(A1) The process $\mu : [0, 1] \rightarrow \mathbb{R}^{p \times q}$ is continuous and $\text{span}(\mu(u), u \in [0, 1]) = E_c$.

(A2) The process $\hat{\mu} : [0, 1] \rightarrow \mathbb{R}^{p \times q}$ is such that

$$n^{1/2}(\hat{\mu} - \mu, \hat{Q}_k - Q_k) \text{ converges weakly in } l^\infty([0, 1])^{(p \times q)} \times \mathbb{R}^{p \times p} \text{ to } (w_\mu, w_Q).$$

(A3) The process $\hat{\mu}^* : [0, 1] \rightarrow \mathbb{R}^{p \times q}$ is such that, conditionally on the sample,

$$n^{1/2}(\hat{\mu}^* - \hat{\mu}, \hat{Q}_k^* - \hat{Q}_k) \text{ converges weakly in } l^\infty([0, 1])^{(p \times q)} \times \mathbb{R}^{p \times p} \text{ to } (\tilde{w}_\mu, \tilde{w}_Q),$$

in probability, with $(Q\tilde{w}_\mu, \tilde{w}_Q\mu) \stackrel{d}{=} (Qw_\mu, w_Q\mu)$.

The previous set of assumptions might be understood as follows. Assumption (A1) is the so called coverage condition that has been used by several authors [7], [23]. This condition is discussed within the SIR and CUME context in Remark 4. Assumptions (A2) and (A3) depends on the test under consideration. When $k = 2$, for SIR and CUME, (A2) (resp. (A3)) is a straightforward consequence of Theorem 3 (resp. Theorem 4); for order 2 moments based methods, it is implied by Corollary 3 (resp. Corollary 4). The reader might refer to the mentioned theorems to obtain conditions that guarantee (A2) and (A3). For $k = 1, 3$, the theorems we just mentioned are not enough to obtain directly (A2) and (A3) because these conditions involve the joint distribution of the process $\hat{\mu}$ with a certain eigenprojector. However they can be ascertained by the additional use of an asymptotic expansion for eigenprojectors e.g. Lemma 4.1 in [26]. Finally, note that Assumption (A3) is weaker than asking for a complete bootstrap, i.e. that, conditionally on the sample, $n^{1/2}(\hat{\mu}^* - \hat{\mu}, \hat{Q}_k^* - \hat{Q}_k)$ has the same weak limit as $n^{1/2}(\hat{\mu} - \mu, \hat{Q}_k - Q_k)$, in probability. This will have interesting consequences on the validity of different bootstrap strategies (see the remark bellow).

Proposition 7. *Under Assumptions (A1), (A2) and (A3), testing (9), (10) or (11) with respectively $\hat{\Lambda}_1, \hat{\Lambda}_2, \hat{\Lambda}_3$ and calculation of the quantiles with $\hat{\Lambda}_1^*, \hat{\Lambda}_2^*, \hat{\Lambda}_3^*$ respectively, is consistent.*

Proof. Note that $\hat{\Lambda}_k$ is a continuous transformation of the process $n^{1/2}\hat{Q}_k\hat{\mu}$. Under H_0 , because (A1) implies that $Q_k\mu = 0$, we have

$$n^{1/2}\hat{Q}_k\hat{\mu} = n^{1/2}Q_k(\hat{\mu} - \mu) + n^{1/2}(\hat{Q}_k - Q_k)\mu + n^{1/2}(\hat{Q}_k - Q_k)(\hat{\mu} - \mu). \quad (13)$$

Using (A2), $\|\hat{\mu} - \mu\|_\infty \rightarrow 0$ in probability, then by Slutsky's Lemma, the last term vanishes asymptotically. Using (A3) and the continuous mapping theorem, the sum of the first two terms in (13) (and so $n^{1/2}\hat{Q}_k\hat{\mu}$) converges weakly in $l^\infty([0, 1])^{(p \times q)}$. As a consequence of the continuous mapping theorem, under H_0 , $\hat{\Lambda}_k$ converges weakly to a real random variable. Under H_1 , it is easy to show that $|Q_k\mu(u)|_2 > 0$ for a certain $u \in [0, 1]$, making $\hat{\Lambda}_k$ going to infinity in probability.

Consequently it is enough to show that the bootstrap statistic (i) has the same behaviour as the statistic under H_0 , and (ii) remains bounded in probability under H_1 . For (i), note that $\hat{\Lambda}_k^*$ is a continuous transformation of the process $n^{1/2}\hat{Q}_k^*\hat{\mu}_k^*$ that can be written as

$$n^{1/2}\hat{Q}_k^*\hat{\mu}_k^* + n^{1/2}(\hat{Q}_k^* - \hat{Q}_k)\hat{\mu}_k^*$$

then using the definition of $\hat{\mu}_k^*$, we get that

$$n^{1/2}\hat{Q}_k^*\hat{\mu}_k^* = n^{1/2}\hat{Q}_k(\hat{\mu}^* - \hat{\mu}) + n^{1/2}(\hat{Q}_k^* - \hat{Q}_k)(I - \hat{Q}_k)\hat{\mu} + n^{1/2}(\hat{Q}_k^* - \hat{Q}_k)(\hat{\mu}^* - \hat{\mu}).$$

The latter term is asymptotically neglectable by (A3), it follows that

$$n^{1/2}\widehat{Q}_k^*\widehat{\mu}_k^* = n^{1/2}Q_k(\widehat{\mu}^* - \widehat{\mu}) + n^{1/2}(\widehat{Q}_k^* - \widehat{Q}_k)(I - Q_k)\mu + o_p(1).$$

Since under H_0 , $(I - Q_k)\mu(u) = \mu(u)$, using (A3) and the continuous mapping theorem is enough to show that conditionally on the sample, $n^{1/2}\widehat{Q}_k^*\widehat{\mu}_k^*$ has the same asymptotic law as $n^{1/2}\widehat{Q}_k\widehat{\mu}_k$, in probability. Then invoking again the continuous mapping theorem provide the same conclusion with $\widehat{\Lambda}_k^*$ and $\widehat{\Lambda}_k$. Under H_1 , in light of the latter representation and by (A3), conditionally on the sample, the sequence $n^{1/2}\widehat{Q}_k^*\widehat{\mu}_k^*$ is tight. \square

Remark 7 (other bootstrap strategies). As we have highlighted (see Remarks 2 and 5), the natural bootstrap candidate for $\widehat{c}_{\widehat{F}}$ is given by $\widehat{c}_{\widehat{F}^*}^*$ (rather than $\widehat{c}_{\widehat{F}}^*$), in which the estimated cdf \widehat{F} has been bootstrapped. Because the randomness of the slices (carried by \widehat{F}) affects the limiting distribution, this can be seen, at first glance, as a necessary evil. In our particular context given by (9), (10) and (11), and under the linearity condition, it is in fact not essential to bootstrap \widehat{F} . Indeed, Assumptions (A3) only requires that the bootstrap estimator reproduces the law of $Q\sqrt{n}(\widehat{c}_{\widehat{F}} - c_F)$ where Q stands for the orthogonal projector on a given subspace of E_c^\perp . In light of the proof of Theorem 3, we have that $\sqrt{n}(\widehat{c}_{\widehat{F}} - c_F)$ has the following limiting distribution

$$w_1 - \partial c_F B,$$

where (w_1, B) is a certain Gaussian process. Since for any $u \in [0, 1]$, $\partial c_F(u) = E(Z|Y = F^-(u))$, using the linearity condition we have that $\partial c_F \in E_c$. Multiplying by Q the latter representation, we obtain that the asymptotic law of $Q\sqrt{n}(\widehat{c}_{\widehat{F}} - c_F)$ is reduced to the representation Qw_1 . As a consequence, the part $\partial c_F B$ in the asymptotic variance does not matter here, and so the bootstrap estimator given by $\widehat{c}_{\widehat{F}}^*$, satisfies assumption (A3) as well as $\widehat{c}_{\widehat{F}^*}^*$ does. Either for SIR or CUME, using $\widehat{c}_{\widehat{F}}^*$ is computationally less intensive than using $\widehat{c}_{\widehat{F}^*}^*$ because it preserves the slicing initially used for the estimator $\widehat{c}_{\widehat{F}}$.

4 Simulations

In this section, we study the accuracy of the bootstrap approximation facing one of the Cramér-von Mises tests introduced in Section 3.3. We focus on the test of significance of some sets of predictors $\eta^T X$ described by (10) and we consider the performance of both methods SIR and CUME with the statistic $\widehat{\Lambda}_2^*$. Our aim is to analyse quite difficult situations from small to moderate sample size.

Given i.i.d. observations from a regression model, we test whether a vector η is orthogonal to E_c or not. The statistics of interest are related to SIR with H slices and CUME, each is given respectively by

$$\begin{aligned}\widehat{\Lambda}_2^{\text{SIR}} &= n \int |\widehat{Q}_\eta \widehat{m}(u, H^{-1})|_F^2 d\nu_d(u) \\ \widehat{\Lambda}_2^{\text{CUME}} &= n \int |\widehat{Q}_\eta \widehat{c}_{\widehat{F}}(u)|_F^2 d\nu_c(u),\end{aligned}$$

where \widehat{Q}_η is the orthogonal projector on the space generated by $\widehat{\Sigma}^{-1/2}\eta$, $\widehat{\Sigma}$ is the classical estimator of the variance of X , and ν_d (resp. ν_c) is the uniform probability measure on the set $\{(2h-1)/(2H), h=1, \dots, H\}$ (resp. on the set $[0, 1]$). The process $\widehat{c}_{\widehat{F}}$ and \widehat{m} are the same as the ones define in the paper except that from now on, we estimate the mean and the variance of X .

The bootstrap estimators are computed following Equation (12). As pointed out in Remark 7, there are two different bootstrap strategies that are available to compute the quantiles of the test. The first one involves $\widehat{c}_{\widehat{F}^*}^*$ and gives, for instance, for CUME

$$n \int |\widehat{Q}_\eta^* \{\widehat{c}_{\widehat{F}^*}^*(u) - \widehat{Q}_\eta \widehat{c}_{\widehat{F}}(u)\}|_F^2 d\nu_c(u),$$

where \widehat{Q}_η^* is the orthogonal projector on the space generated by $\widehat{\Sigma}^{*-1/2}\eta$ and $\widehat{\Sigma}^*$ is defined in Remark 9 bellow. This bootstrap is abbreviated in the next SIRb1 and CUMEb1. The second bootstrap involves $\widehat{c}_{\widehat{F}}^*$ and gives, for instance, for CUME

$$n \int |\widehat{Q}_\eta^* \{\widehat{c}_{\widehat{F}}^*(u) - \widehat{Q}_\eta \widehat{c}_{\widehat{F}}(u)\}|_F^2 d\nu_c(u),$$

it is abbreviated by SIRb2 and CUMEb2. To compute a quantile of level α , we draw independently B bootstrap statistics and then calculate the empirical quantile of level α associated to this sample.

Remark 8 (computation of CUME). Either for the estimator or the bootstrap, integrals associated to CUME are computed easily because the integrands are piecewise constant. For instance, one may show that $\widehat{\Lambda}_2^{\text{CUME}} = n^{-1} \sum_{i=1}^n |\widehat{Q}_\eta \widehat{c}_{id}(Y_i)|_2^2$, and the same kind of formulas can be derived for the bootstrap statistics. Because, the integrand of the method b1 has $2n$ jumps whereas the integrand of the method b2 has n jumps, the method b2 is less intensive computationally.

Remark 9 (standardizing the bootstrap). For the sake of completeness, in this section we have leaved the theoretical framework of the paper that supposed to be known the mean and the variance of X . To build our estimators, we have plugged the classical estimators of the latter quantities in the initial estimators. This naturally induces an additional part in the asymptotic distribution. We account for this part by bootstrapping also the mean and the variance by respectively

$$\overline{wX} = n^{-1} \sum_{i=1}^n w_{i,n} X_i \quad \text{and} \quad \widehat{\Sigma}^* = n^{-1} \sum_{i=1}^n w_{i,n} (X_i - \overline{wX})(X_i - \overline{wX})^T.$$

Note that they are used to standardized the predictors as well as to standardized the set of directions η under test.

We consider the following models:

$$Y = X_1 + \sigma e, \tag{14}$$

$$Y = \frac{X_1}{0.5 + (2 + X_2 + X_3)^2} + \sigma e, \tag{15}$$

$$Y = \exp(X_1) \times \sigma e, \tag{16}$$

where $(X, e) \in \mathbb{R}^5$ follows a standard normal distribution. Model (14) has already been considered in Section 3.1 in order to highlight the influence of the randomness of the slices on the asymptotic distribution of the estimators. Model (15) is borrowed from [19] and Model (16) represents a regression model with non-additive noise. Variations of σ permits to switch from easy to more difficult situations. We have ran 1000 Monte-Carlo replication, for which we have performed the test under H_0 (when $\eta = (0, 0, 0, 1)$) and under H_1 (when $\eta = (1, 0, 0, 0)$) with SIRb1&2 and CUMEb1&2 at the nominal level of $\alpha = 5\%$. In each case, the bootstrap sample

σ	n	H_0	SIRb1				CUMeb1	SIRb2				CUMeb2
			(H) 3	5	7	10		(H) 3	5	7	10	
.5	30	T	14	3	0	0	20	156	192	207	224	173
		F	1000	999	988	804	1000	1000	1000	1000	1000	1000
	50	T	15	2	0	0	27	107	129	123	116	108
		F	1000	1000	1000	1000	1000	1000	1000	1000	1000	1000
	100	T	21	7	3	0	25	71	107	79	85	81
		F	1000	1000	1000	1000	1000	1000	1000	1000	1000	1000
1	200	T	39	14	8	0	31	72	69	66	63	57
		F	1000	1000	1000	1000	1000	1000	1000	1000	1000	1000
	30	T	17	0	0	0	38	151	180	199	213	151
		F	935	794	591	186	990	978	962	960	921	993
	50	T	31	1	0	0	61	132	139	141	136	128
		F	1000	991	971	809	1000	1000	999	1000	994	1000
.5	100	T	22	8	2	0	54	69	64	90	87	84
		F	1000	1000	1000	1000	1000	1000	1000	1000	1000	1000
	200	T	34	16	4	1	60	77	82	67	76	89
		F	1000	1000	1000	1000	1000	1000	1000	1000	1000	1000

Table 1: Estimated level and power in Model (14) for $\alpha = 5\%$ with 1000 replications.

σ	n	H_0	SIRb1				CUMeb1	SIRb2				CUMeb2
			(H) 3	5	7	10		(H) 3	5	7	10	
.1	30	T	11	2	0	0	36	175	191	209	213	174
		F	993	946	870	450	1000	1000	997	995	995	1000
	50	T	17	1	0	0	33	125	111	99	131	121
		F	1000	1000	1000	983	1000	1000	1000	1000	1000	1000
	100	T	26	4	0	0	32	90	75	100	81	86
		F	1000	1000	1000	1000	1000	1000	1000	1000	1000	1000
.5	200	T	30	21	3	0	39	63	67	71	69	67
		F	1000	1000	1000	1000	1000	1000	1000	1000	1000	1000
	30	T	16	1	0	0	76	156	162	188	197	161
		F	553	296	137	17	786	736	732	723	661	835
	50	T	27	2	0	0	63	115	114	133	107	103
		F	803	672	498	176	946	897	892	872	832	956
.5	100	T	30	11	2	1	72	79	81	109	83	96
		F	989	977	954	821	1000	998	998	993	991	1000
	200	T	32	17	3	0	48	65	59	63	72	62
		F	1000	1000	1000	1000	1000	1000	1000	1000	1000	1000

Table 2: Esimtated level and power in Model (15) for $\alpha = 5\%$ with 1000 replications.

σ	n	H_0	SIRb1				CUMeb1	SIRb2				CUMeb2
			(H) 3	5	7	10		(H) 3	5	7	10	
.5	30	T	23	3	0	0	80	148	190	208	214	134
		F	511	410	241	25	264	847	940	946	900	780
	50	T	23	0	0	0	72	115	130	134	124	107
		F	895	924	857	553	748	974	993	993	993	966
	100	T	35	5	1	1	66	79	80	89	79	90
		F	997	1000	1000	1000	1000	999	1000	1000	1000	1000
	200	T	35	21	5	1	60	67	63	66	67	70
		F	1000	1000	1000	1000	1000	1000	1000	1000	1000	1000
1	30	T	27	3	0	0	93	159	179	190	236	158
		F	518	441	277	29	289	834	926	939	912	794
	50	T	23	3	2	0	71	100	99	111	119	96
		F	897	935	865	547	739	974	997	996	996	974
	100	T	27	8	0	0	65	85	84	101	94	81
		F	999	1000	1000	1000	1000	1000	1000	1000	1000	1000
	200	T	37	14	3	0	52	75	61	69	72	61
		F	1000	1000	1000	1000	1000	1000	1000	1000	1000	1000

Table 3: Esimtated level and power in Model (16) for $\alpha = 5\%$ with 1000 replications.

number was equal to $B = 500$. The number of rejections in each situation is given along the tables 1 to 3.

The conclusions might be drawn as follows, first comparing SIR and CUME and second evaluating the differences between b1 and b2.

Comparing SIR and CUME, we must raise from the start, that contrary to SIR, CUME no longer depends on the number of slices H . Unfortunately, we see that SIR is strongly affected by changes in H , notably at small sample sizes. For instance, under the null in Model (16), as soon as H is large, b1 no longer rejects the null while b2 reject the null 20% of the times. Looking at the complete picture offered by all the tables, the smaller H the better, making the SIR 3-slices approach the best competitor for facing CUME. In spite of this “a posteriori” and advantageous selection of H , SIR does not perform better than CUME. For additive models (14) and (15), it seems preferable to use CUME over SIR whereas for Model (16), the situation is slightly mitigated by the high power provided by SIR. To conclude, CUME offers a more simple (no selection of H) approach than SIR and among the considered models, it is more accurate to test with CUME rather than SIR.

Both bootstrap tests b1 and b2 converge to the nominal level but with (in average) opposite signs, i.e. b1 tends to underestimate the level while b2’s estimated level is always greater than the nominal one. This suggests that b1 is more conservative than b2. Meanwhile, the power associated to b2 is always greater than the power associated to b1. Then and in particular for SIR, both are difficult to compare, and the choice between b1 and b2 should be done with care by the user given the trade-off between conservativeness and powerfulness. For CUME, the situation is rather different than for SIR and the clear winner is b1, notably because of the too high level of type I error committed by b2.

5 Conclusion

We have provided a new approach for inverse regression based on empirical processes. This approach has offered a precise description of the asymptotic behaviour of the estimators as well as the validity of the bootstrap. The framework we develop in the paper is linked with the

class of indicator functions. This choice was convenient since the metric entropy properties of this class are widely known, but also because of the natural link it induced with the popular methods SIR and CUME. However, the approach developed in this paper can be extended to different classes of functions than indicators. Indeed, for the order 1 moments based method, one can consider the vector

$$E[X\psi(Y)],$$

when ψ varies among a certain family of functions. Another subject of interest for further studies is right-censored data. Suppose we observe

$$\min(Y, C) \quad \text{and} \quad \mathbb{1}_{\{Y \leq C\}}, \quad \text{where} \quad Y \perp\!\!\!\perp C | X,$$

variations of SIR have been studied for instance in [20] and [22]. It requires a smoothing procedure in order to take into account the effect of the censoring.

Acknowledgement. The author would like to thank Bernard Delyon for helpful comments and advices on this article. He also thank Zhenghui Feng for sharing the Matlab code of CUME.

References

- [1] M. P. Barrios and S. Velilla. A bootstrap method for assessing the dimension of a general regression problem. *Statist. Probab. Lett.*, 77(3):247–255, 2007.
- [2] M. P. Bentler and J. Xie. Corrections to test statistics in principal hessian directions. *Statist. Probab. Lett.*, 47(4):381–389, 2000.
- [3] Caroline Bernard-Michel, Laurent Gardes, and Stéphane Girard. Gaussian regularized sliced inverse regression. *Stat. Comput.*, 19(1):85–98, 2009.
- [4] E. Bura and J. Yang. Dimension estimation in sufficient dimension reduction: a unifying approach. *J. Multivariate Anal.*, 102(1):130–142, 2011.
- [5] R. D. Cook. *Regression graphics*. Wiley Series in Probability and Statistics: Probability and Statistics. John Wiley & Sons Inc., New York, 1998.
- [6] R. D. Cook and B. Li. Dimension reduction for conditional mean in regression. *Ann. Statist.*, 30(2):455–474, 2002.
- [7] R. D. Cook and L. Ni. Sufficient dimension reduction via inverse regression: a minimum discrepancy approach. *J. Amer. Statist. Assoc.*, 100(470):410–428, 2005.
- [8] R. D. Cook and S. Weisberg. Discussion of “sliced inverse regression for dimension reduction”. *J. Amer. Statist. Assoc.*, pages 28–33, 1991.
- [9] R. Dennis Cook. Testing predictor contributions in sufficient dimension reduction. *Ann. Statist.*, 32(3):1062–1092, 2004.
- [10] B. Efron. Bootstrap methods: another look at the jackknife. *Ann. Statist.*, 7(1):1–26, 1979.
- [11] J. Fermanian, D. Radulovic, and M. Wegkamp. Weak convergence of empirical copula processes. *Bernoulli*, 10(5):847–860, 2004.

- [12] P. Hall and K. Li. On almost linearity of low-dimensional projections from high-dimensional data. *Ann. Statist.*, 21(2):867–889, 1993.
- [13] P. Hall and S. R. Wilson. Two guidelines for bootstrap hypothesis testing. *Biometrics*, 47(2):757–762, 1991.
- [14] W. Härdle and T. M. Stoker. Investigating smooth multiple regression by the method of average derivatives. *J. Amer. Statist. Assoc.*, 84(408):986–995, 1989.
- [15] J. Hoffmann-Jorgensen. *Stochastic processes on Polish spaces*, volume 39 of *Various Publications Series (Aarhus)*. Aarhus Universitet, Matematisk Institut, Aarhus, 1991.
- [16] M. Hristache, A. Juditsky, and V. Spokoiny. Direct estimation of the index coefficient in a single-index model. *Ann. Statist.*, 29(3):595–623, 2001.
- [17] T. Hsing and R. J. Carroll. An asymptotic theory for sliced inverse regression. *The Annals of Statistics*, 20(2):1040–1061, 1992.
- [18] B. Li and S. Wang. On directional regression for dimension reduction. *J. Amer. Statist. Assoc.*, 102(479):997–1008, 2007.
- [19] K. Li. Sliced inverse regression for dimension reduction. *J. Amer. Statist. Assoc.*, 86(414):316–342, 1991.
- [20] K. Li, J. Wang, and C. Chen. Dimension reduction for censored regression data. *Ann. Statist.*, 27(1):1–23, 1999.
- [21] Y. Li and L. Zhu. Asymptotics for sliced average variance estimation. *Ann. Statist.*, 35(1):41–69, 2007.
- [22] N. V. Nadkarni, Y. Zhao, and M. R. Kosorok. Inverse regression estimation for censored data. *Journal of the American Statistical Association*, 106(493), 2011.
- [23] F. Portier and B. Delyon. Optimal transformation: a new approach for covering the central subspace. *J. Multivariate Anal.*, 115:84–107, 2013.
- [24] F. Portier and B. Delyon. Bootstrap Testing of the Rank of a Matrix via Least-Squared Constrained Estimation. *J. Amer. Statist. Assoc.*, 109(505):160–172, 2014.
- [25] J. Præstgaard and J. A. Wellner. Exchangeably weighted bootstraps of the general empirical process. *Ann. Probab.*, 21(4):2053–2086, 1993.
- [26] D. E. Tyler. Asymptotic inference for eigenvectors. *Ann. Statist.*, 9(4):725–736, 1981.
- [27] A. W. van der Vaart and J. A. Wellner. *Weak convergence and empirical processes*. Springer Series in Statistics. Springer-Verlag, New York, 1996. With applications to statistics.
- [28] X. Yin and R. D. Cook. Dimension reduction for the conditional k th moment in regression. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 64(2):159–175, 2002.
- [29] L. Zhu and K. Fang. Asymptotics for kernel estimate of sliced inverse regression. *Ann. Statist.*, 24(3):1053–1068, 1996.
- [30] L. Zhu and K. W. Ng. Asymptotics of sliced inverse regression. *Statist. Sinica*, 5(2):727–736, 1995.
- [31] L. Zhu, L. Zhu, and Z. Feng. Dimension reduction in regressions through cumulative slicing estimation. *J. Amer. Statist. Assoc.*, 105(492):1455–1466, 2010.